

Cleansing procedures for overlaps and inconsistencies in administrative data: the case of German labour market data

Oberschachtsiek, Dirk; Scioch, Patrycja

Veröffentlichungsversion / Published Version

Zeitschriftenartikel / journal article

Zur Verfügung gestellt in Kooperation mit / provided in cooperation with:

GESIS - Leibniz-Institut für Sozialwissenschaften

Empfohlene Zitierung / Suggested Citation:

Oberschachtsiek, D., & Scioch, P. (2009). Cleansing procedures for overlaps and inconsistencies in administrative data: the case of German labour market data. *Historical Social Research*, 34(3), 242-259. <https://doi.org/10.12759/hsr.34.2009.3.242-259>

Nutzungsbedingungen:

Dieser Text wird unter einer CC BY Lizenz (Namensnennung) zur Verfügung gestellt. Nähere Auskünfte zu den CC-Lizenzen finden Sie hier:
<https://creativecommons.org/licenses/by/4.0/deed.de>

Terms of use:

This document is made available under a CC BY Licence (Attribution). For more Information see:
<https://creativecommons.org/licenses/by/4.0>

Cleansing Procedures for Overlaps and Inconsistencies in Administrative Data. The Case of German Labour Market Data

*Patrycja Scioch & Dirk Oberschachtsiek**

Abstract: »Die Bedeutung von Bereinigungsverfahren für Überschneidungen und Inkonsistenzen in administrativen Daten. Am Beispiel deutscher Arbeitsmarktdaten«. Process-generated and administrative datasets have become increasingly important for labour market research over the past ten years. Major advantages of this data are large sample sizes, absence of retrospective gaps and unit nonresponses. Nevertheless, the quality and validity of the information remain unclear. This paper contributes to this subject, focusing on the variation of research results due to alternative data cleansing procedures. In particular, the paper uses the general set up for data cleaning proposed by Wunsch/Lechner (2008) in evaluating the outcome of training programmes in Germany. First results are limited to the sensitivity of the construction of the sample populations used for the counterfactuals analysis. The results emphasize that sample construction seems to be robust to the scenario used for the data cleansing.

Keywords: Longitudinal Analysis, Process-Generated Data, Social Bookkeeping Data, Public Administrative Data, Data Management Record Linkage, Data Fusion, Labour Market Data.

1. Motivation

Process-generated and administrative datasets have become increasingly important in research over the past ten years. Kluve (2006), for example, reports that almost 80% of all microeconomic evaluation studies in Europe are based on that type of data. Particularly Scandinavian labour market research is based on register data (Eliason and Storrie 2006; Carling and Richardson 2004; Roed and Raaum 2003).

Administrative data have various advantages; besides providing extensive information on individuals, administrative data can help overcome some weak-

* Address all communications to: Patrycja Scioch, Institute for Employment Research (Research Data Centre), Regensburger Str. 104, 90478 Nuremberg, Germany; e-Mail: patrycja.scioch@iab.de.

Dirk Oberschachtsiek, Social Science Research Centre Berlin (WZB), Reichpietschufer 50, 10785 Berlin, Germany ; e-Mail: Dirk.Oberschachtsiek@wzb.eu.

We thank Conny Wunsch and Michael Lechner for supporting the work related to this paper and giving access to their programme codes. In addition we thank Nina Baur for helpful comments.

nesses of survey data like attrition bias, reporting or recollection bias, the lack of relevant comparison groups and small sample size. One of the most important advantages of administrative data concerns the option of merging and combining information from different administrative data sources¹ and over multiple points of time. However, the quality of information might suffer as a result of merging. It depends crucially on the consistency of identifiers and information coming from the different sources.

So far there are only few studies that focus on the quality of administrative data. With respect to survey data this has been the subject of research for 20 years (e.g. Schnell, 1985, 1991). With regard to process-produced data first analyses concentrate on the data generating process and its complexity (Kruppe and Oertel 2003; Engelhardt et al. 2008). Further studies show that there are similar problems like missing values, overlaps and inconsistencies. Jaenichen et al. (2005) or Bernhard et al. (2006) refer to the requirement of data preparation and data cleansing. Recent work also focuses on the connection of research results and data cleansing procedures (e.g. Kruppe et al. 2008; Waller 2007).

This study contributes to the latter type of research. Based on the Integrated Employment Biographies (IEB – Integrierte Erwerbsbiographien) we investigate the impact of different cleansing procedures on data overlaps and inconsistencies between different data sources. The IEB data are compiled from four distinct and independent administrative sources stemming from the German Employment Services and have been used quite extensively for the evaluation of active labour market policies (ALMPs) in Germany (e.g. Biewen et al. 2007; Wunsch and Lechner 2008). In order to analyse the effect of different data cleansing procedures we use one of these evaluation studies and replicate the results based on different variations of the data cleansing methods suggested by Wunsch/Lechner (2008).

This paper is structured as follows. In the next section we describe the database and discuss problems that may occur when using the data. Section 3 presents previous studies on identifying and handling of inconsistencies and overlaps in the IEB before section 4 explains the replication and identification of variance. Section 5 summarises the descriptive results to value the quality of replication and the variance in the evaluation samples. Section 6 concludes.

2. The Database: Integrated Employment Biographies (IEB)

The database used in this study is the Integrated Employment Biographies (IEB) of the Institute for Employment Research (IAB), which is a longitudinal data set merged from four distinct process-generated data sources. The data

¹ Merging is possible with identifiers on individual level.

cover nearly 80% of the total labour force in Germany and almost 100% of the employees liable to social security. Not included are periods of self-employment, civil servants and periods of childcare leave. For detailed information see Jacobebbinghaus/Seth (2007).

The data set's four sources are fed by four administrative processes, and linked together by using a unique identifier that allows to combine the observations. Each of these sources offers a brought set of attributes and covers different periods of observation.

- 1) The first data source is the Employment Histories (*Beschäftigten-Historik*) containing employment periods captured by the social insurance register back until 1990. Beside begin and end dates it also includes the employment state, personal characteristics, wage, type of profession, region and the industry. Moreover it allows merging further employer information.
- 2) The second data source contains data on spells of unemployment from the Benefit-Recipient-History (*Leistungsempfänger-Historik*). It has information, on a daily basis, on unemployment benefits, unemployment assistance and subsistence allowances since 1990. Additionally, the source includes personal characteristics and statements on sanctions.
- 3) Most of the individual characteristics in the IEB data arise from the Applicants-Pool data (*Bewerberangebot*), which contains information on job-searching spells since 1999. Apart from the current marital state, nationality, health, education and regional characteristics the data set also comprises information about the last job and on the desired job and profession.
- 4) Finally the data set on active labour market programmes participation (*I-SAAK – Instrumente Aktiver Arbeitsmarktpolitik* or *MTH – Maßnahme-Teilnahme Historik*) provides information on periods spent in promoted schemes. Since 2000 any participation in employment or training measures has been recorded (begin and end date and the characteristics of respective participants and programme).²

The IEB data are organised on a daily basis and allow to control for time varying covariates. Due to the huge size of the IEB the Institute for Employment Research offers access to a 2.2% random sample called IEBS (Integrated Biographies Sample).³ It is important to note that the sources are not cross-validated, which may cause the existence of parallel observations (overlaps). Individuals can have several jobs at the same time or they might be employed and searching for a new job or receiving benefits while on job search or participating in labour market programmes. These spells can be completely parallel, one may embed the other or they are overlapping.

² For a detailed description of the data generating process of the participation in measure data see Engelhardt et al. (2008).

³ See for the data access <http://fdz.iab.de/> and description Zimmermann et al. (2007).

The existence of parallel observations is twofold: It may offer additional information, like periods of promoted employment (see Huber/Schmucker in this Special Issue). However, it may also cause problems when information is contradictory. In the latter case one must decide which data source to believe – which is the subject of data cleansing procedures.

3. Previous Work on Identifying and Handling Inconsistencies in the IEB

There is already a small body of research on identifying inconsistencies in the IEB. Some of the studies investigate them in general and suggest only simple options to deal with overlaps and contradictions whereas the latest ones concentrate on special problems: One of the first studies that address inconsistencies in the IEB is Jaenichen et al. (2005). In a simple framework it tries to identify distinctive types of implausible cases and discusses simple heuristics to handle these types of inconsistencies. In general, the paper focuses on overlaps, gaps and the missing of parallel observations between two of the four sources respectively or within one source. In a second step they draw a subsample of 30 to 50 individuals for each type of implausibility in order to gain potential interpretations and explanations. As a global heuristic they do not find any convenient and robust rule that fits to the variety of inconsistencies under investigation and recommend the application of project specific approaches.

Bernhard et al. (2006) extend this by a comprehensive investigation of all possible overlaps in information in the IEBS⁴. They give an overview over the most common overlaps within and between the sources and define overlapping-types. By means of some examples they discuss possible causes and ways to deal with the contradictions. Furthermore they analyse inconsistencies between two special sources, the Employment History and the Benefit-Recipient History, in more detail by using additional information. They provide information to other researchers to evaluate the meaning and legitimacy of overlaps.

In contrast to these studies which address inconsistencies in general, Kruppe et al. (2008) focus on a single variable and its variance based on different legal definitions, administrative procedures and the validity of the information. They examine six different implementation strategies for common definitions of unemployment in the IEBS. These concepts yield 63 definitions with huge differences in the mean unemployment duration varying between 127 and 325 days of unemployment. They adduce these differences as evidence that the underlying concept of unemployment definition is crucial for applied research.

⁴ IEBS V 1.0 based on the IEB V 3.

Likewise Waller (2007) also concentrates on a single variable in the IEB⁵. Contrary, her focus is the variance of the end dates in program participation on the estimation of treatment effects due to measurement errors. She develops four different correction procedures and discusses the influence on estimation results using different methods (descriptive attendance and employment rates, statistical matching, descriptive duration method). Waller (2007) finds only little differences in the treatment effects caused by measurement errors in the end dates. Significant effects are limited to the lock-in periods and in particular found for long programs. This emphasizes to put effort into the correction of the end dates only if the interest is concerned with exact magnitude of lock-in-effects.

We use this existing body of research as a starting point in order to ask: How big is the effect of cleansing procedures on later statistical results? How complicated do cleansing have to be in order to provide “good” data?

4. The Cleansing Procedures

4.1 The Procedure Suggested by Wunsch and Lechner

The general framework for the data cleansing procedure used in this study has been proposed by Wunsch/Lechner (2008) which is a study on evaluating training and employment programmes to assess the effectiveness of labour market programmes in West Germany. They perform matched pairs comparisons. This procedure allows a simple identification of counterfactual observations since it uses statistical twins with respect to the likelihood to participate in a certain promotion scheme (see Rosenbaum and Rubin 19985). However, this technique needs strong assumptions about relevant characteristics that affect selection and potential outcomes (for details and a deeper discussion see Heckman et al. 1998; Imbens 2004; Caliendo and Kopeinig 2006).

Wunsch/Lechner use a 2% random sample of the IEB supplemented with additional characteristics taken from the different data sources as well as characteristics from regional statistics. The final data set contains personal characteristics and spell related information. For a detailed description of the data see table A1 (appendix).

Based on this data they identify potential comparisons between participants and non-participants taken from the total of inflows into unemployment between January 2000 and December 2002. In order to reduce heterogeneity they restrict the sample to individuals that received unemployment benefits. Participants are limited to individuals who have started a programme during the next 18 months after becoming unemployed and who are receiving unemployment

⁵ IEB V 2.05.

benefits directly before starting the programme. To identify the potential outcome as a difference of spell lengths between treated and the matched non-treated they impute a reference date (the non-observed begin date) by using regression methods.

However, this identification set up still needs to clearly identify one state at each point in time. Wunsch/Lechner (2008) define time frames and rules of priority for possible parallel states. Afterwards they transform the data into a panel data set with exact one state at each point of time. Labour market programmes are treated with the highest priority (followed by periods of benefit receipt and times of employment). The lowest priority gets information out of applicants' pool data.

4.2 Applying the Wunsch/Lechner-Procedure

Referring to Wunsch/Lechner (2008) we use the same set up to construct a matched comparisons analysis based on a more recent draw of the IEBS. Most importantly we use the proposed framework to produce multiple subsamples based on different rules of priorities for individuals with overlapping observations. This results in different final states for the individuals and thus causes variation in the composition of the subsamples used for the matched comparisons analysis.

Similar to Wunsch/Lechner we organise the data in a panel set-up by splitting the spell data into frames of two weeks.⁶ Within these time frames it is now possible to isolate one state:

- 1) *Sorting Rule 1 (Length Priority)*: First all parallel observations are sorted by length.
- 2) *Sorting Rule 2 (Source Priority)*: If two or more parallel observations have the same length we use the respective data source as a proxy of the validity to order the observations.

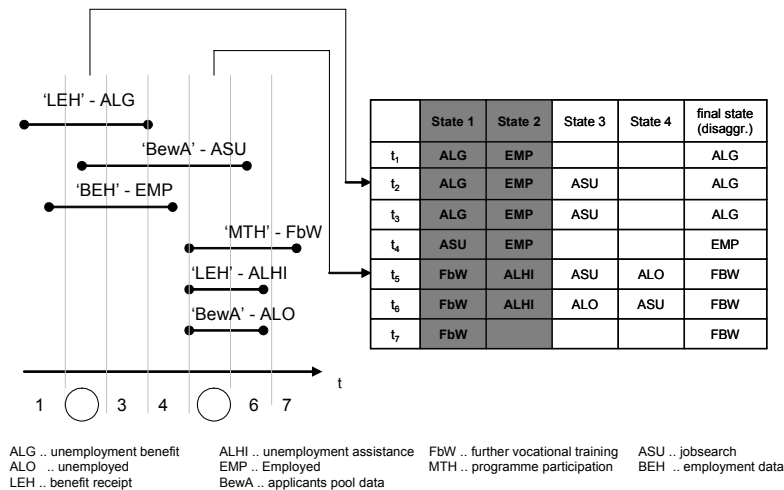
An illustration of this approach is given in figure one, where the left part displays overlapping observations from different sources. The table on the right hand side in figure one shows the data matrix related to this example. As noted above, the period is divided into (seven) time windows. For the whole period we observe six different observations each coming from a distinct source. The first abbreviation specifies the data source and the second one the concrete state. For example: There is a LEH-spell going from time-frame one until the end of time-frame three and a BEH-spell beginning in time-frame one and ending in the mid of time-frame four. The focus of the data cleansing is the identification of one valid observation per time-frame. The right hand side shows the way the observations are transformed into a data matrix, with one

⁶ Wunsch/Lechner made some sensitivity analysis trying different kinds of time windows – shorter and longer than two weeks – without significant differences.

row representing one time-frame and each state in one column (e.g. time-frame one – see t_1 – covers two states and period five contains four states – see t_5 –).

One of the most important steps in this approach refers to the sorting routines. Note that the order of state across the column displayed in figure one is crucial. The first row displays two observations, one coming from the receipt of benefit source and the other from the employment histories. The first column contains the observation with the longest period in window one. If we observe multiple observations with the same length we need to sort the observations by heuristic routines (see time-frame five). Since we are interested in the evaluation of training schemes we may classify all spells with participation in a training scheme with the highest priority. Observations coming from the job search register (note there are two possible states: searching and unemployed) are less valid because they are not associated with any type of payment and are therefore classified with a lower priority. The priority rules suggested by Wunsch /Lechner (2008) are: 1. programme (MTH), 2. benefit receipt (LEH), 3. employment (BEH), 4. job search (BewA).

Fig.1: Identification of the final state



For any further data cleansing we account for the first two states – independent of the number of overlaps. Both selected states are now sorted only based on the Sorting Rule 2 (source priority). To demonstrate the choice of the final state the example continues in the illustration. In time window two we observe an observation coming from the unemployment benefit register and an employment episode. Following the rule of source priority we define the first of these episodes as the final state. Likewise, in period five the final state (fur-

ther vocational training) arises because unemployment assistance has a lower priority than the participation in a labour market programme. Results are displayed in column five of figure one.

4.3 Variation of Sorting and Priority Rules

It is obvious, that the sorting and priority rules influence the status a person is assigned to in the final data set. Changing the priority of the data sources alters the definitions of the final states which leads to different data samples. In the rest of this paper, we want to test how strong the influence of these cleansing procedures is on the final data set. In order to do so, we compare three different methods of data cleansing procedures:

- 1) *Method V0: 1. programme (MTH), 2. benefit receipt (LEH), 3. employment (BEH), 4. job search register data (BewA).*

Method V0 follows the order of priority of Wunsch/Lechner (2008). As mentioned above, when evaluating labour market programmes the participation in a programme should get the highest priority. Sources associated with payments (benefit-recipient-history (LEH) and the employment history (BEH)) follow on second and third priority. By contrast, the job search register contains a lot of optional information and is considered to be less valid. This leads to the last priority of episodes coming from this data source. Method V0 is used as the reference method.

- 2) *Method V1: 1. programme (MTH), 2. employment (BEH), 3. benefit receipt (LEH), 4. job search register data (BewA).*

The first variation occurs in method V1, where the priority of the two sources with money payments is reversed. Both are regarded as valid and there is no clear indication which one to prefer. Altering the priority of both sources may lead to a significant change of the number and duration of employment spells in the analysis sample.

- 3) *Method V2: 1. employment (BEH), 2. programme (MTH), 3. benefit receipt (LEH), 4. job search register data (BewA).*

Method V2 assumes that the participants-in-measure database is not considered to be fully valid. To some extent all dates of this source may be considered to be planned data. Usually, the information related to participations is collected when the programme is assigned to the individual. If the programme is cancelled, delayed or the individual does not take part this is not updated in every case. Thus, method two degrades the priority of this data. However, since some participations come along with benefits and the interest of any evaluation focuses on the effects of participation we order participation as priority two. Assigning participations behind LEH would lead to a dramatic reduction of the participations used for subsequent evaluation studies.

5. Results:

Effect of Sorting Rules on the Final Data Set

The results described below are restricted to descriptive findings focusing on the difference between the composition of the evaluation data samples. The first part of this section presents the quality of the replication of Wunsch and Lechner's (2008) results by comparing means statistics. The second part focuses on the variance in the different evaluation samples (V0 vs. V1; V0 vs. V2).

5.1 Quality of the Replication

The goodness of the replication can be assessed by comparing the means and shares of the variables of the Wunsch/Lechner evaluation sample.⁷ Unfortunately, the means are listed without decimal places and standard deviation. For simplicity we assume the variance to be the same in the original and the replicated data (sample V0). Given this variance we calculate the confidence interval of the original data which can be used for a rough comparison of both samples.

Compared to the original data used in Wunsch/Lechner (2008) the number of observations is higher in sample V0. Only for general further training with duration below six months (≤ 6 months) the number of cases decreases to half of the amount in the original data. The increase can be explained with a more recent version of the underlying data set, but until now we could not find any reasons for the decrease in only one programme type.

Moreover we tested 140 variables from personal characteristics like age, gender, nation, marital status, disability, health problems, education, apprenticeship and related subcategories, as well as information about the desired job, the profession, status and earnings in last job, the remaining unemployment benefit claim, characteristics of the employment history over the 10 years before entering unemployment and a wide range of regional information.

The number of significant differences varies between 39 (27.9%) in degree courses and 73 (52.1%) for general further training (≤ 6 months). Most of the significant differences occur in the employment history variables and the regional information. Note that some of the variables relate to each other, e.g. the occupational sector of the desired job for example is parted into six subcategories. This inter-correlation leads to an overestimation of the sample differences when focusing on the number of variables with significant differences.

The distribution of the significant differences in dummy-variables is displayed in figure A1 (see appendix). The figure shows the frequency distribution

⁷ In method V0 the approach of Wunsch/Lechner was completely replicated in terms of data preparation and particularly in the order of priority.

of the significant differences and the kernel density estimates. The main part of the differences remains below a value of 0.1 indicating that they do not exceed 10 percentage points. Most of them relate to variables which display the share of individuals in different groups like „local unemployment rate is below 5%“, „between 5% and 7.5%“, etc. As mentioned above, a difference in one of these subcategories leads to significant differences in another (related) variable. Overall, the results indicate differences in the samples, but the low magnitude can be interpreted as a satisfying approximation to the sample used in Wunsch/Lechner (2008).⁸

5.2 Variance in the Evaluation Samples

To gauge the influence of the different cleansing procedures we compare the evaluation samples with the different underlying order of priority. Again, the comparison focuses on testing differences of the sample means. We limit the discussion to a selected amount of 14 variables, which are displayed in table A2 for each type of programme⁹ and the group of non-participants.¹⁰ The means comparison tests between the different methods are carried out for each programme separately. An analysis of differences between the types of programmes (A-F in table A2) within each method is not the subject of this work. For reasons of simplicity the results are discussed by using short training as example for the others. Only if there are variances of relevance in the other programme types, these cases are discussed separately.

Table A2 shows the means and shares for participants in short training. Each of the three methods is described by the number of observations (n), the mean and the standard deviation. Method V0 – the replication – is used as the reference to methods V1 and V2.

The results show that the number of observations does not vary in a wide range, it differs between the methods but not to a significant extent. This is the same for all characteristics, they differ by values of one to three percentage points or are completely equal.

Summing up, we find no significant differences between the sorting methods, which indicates that simple decision rules are sufficient. Even in cases with significant differences, which occur in the group of non-participants (see table A2), the differences range between one and four percentage points and are not of relevance. For example the duration of last unemployment in method V0 with 4.97 months decreases to 4.71 months in method V1. According to the

⁸ Some variables are far from the original ones and have to be investigated again but the main part of the data has no significant differences.

⁹ Short training, short combined measures, job-related training, jobseeker assessment, general further training ≤ 6 months, general further training > 6 months, degree courses).

¹⁰ Please contact the author for a detailed list of all variables.

test this is a significant difference, but the magnitude is just eight days, so that this difference seems to be of no practical importance.

6. Summary and Discussion

This paper investigates the influence of variations in cleansing procedures on overlaps in a merged administrative data set. The study presents the cleansing methods and the effects of data cleansing that yield to distinctive analysis samples. First, we replicate the preparation procedures applied in Wunsch/Lechner (2008). Despite some differences, which are restricted to certain kinds of variables, the quality of the replication is satisfying. Therefore the resulting sample can be considered as a sufficient approximation of the original data.

In a second step we develop and apply two variations of the cleansing procedures: by changing the order of priority in cases of overlapping observations the decision rule changes and thus also the final states. Thereafter, we study the influence of these different procedures on the resulting samples using mean comparison tests. These tests show that there are no remarkable significant differences, which is consistent with the findings of previous studies (e.g. Waller 2007). But differences occur when providing a basis for the induction of variation (replicating the approach of Wunsch/Lechner). This leads, on the one hand, to the conclusion that the results are relatively robust to variations in data cleansing procedures. On the other hand, there are sources of variance when constructing the sample which are not detected until now.

The implications on point estimators remain unclear and have to be investigated in further studies. Another interesting extension would be the application of a 'naïve' procedure which prefers observations of one distinct data source without considering aspects of possible and allowed combinations of states.

References

- Bernhard, Sarah / Dressel, Christian / Fitzenberger, Bernd / Schnitzlein, Daniel / Stephan, Gesine (2006): Überschneidungen in der IEBS: Deskriptive Auswertung und Interpretation. FDZ Methodenreport 04/2006.
- Biewen, Martin / Fitzenberger, Bernd / Osikominu, Aderonke / Waller, Marie (2007): Which Program for Whom? Evidence on the Effectiveness of Public Sponsored Training Programs in Germany. IZA Discussion Paper 2885.
- Caliendo, Marco / Kopeinig, Silke (2008): Some Practical Guidance for the Implementation of Propensity Score Matching. In: Journal of Economic Surveys 22 (1). 31-72.
- Carling, Kenneth / Richardson, Katarina (2004): The relative efficiency of labor market programs: Swedish experience from the 1990s. In: Labour Economics 11 (3). 335-354.

- Eliason, Marcus / Storrie, Donald (2006): Lasting or Latent Scars? Swedish Evidence on Long-Term Effects of Job Displacement. In: *Journal of Labor Economics* 24 (4). 831-856.
- Engelhardt, Astrid / Oberschachtsiek, Dirk / Scioch, Patrycja (2008): Datengenese zweier Datenkonzepte: MTG (Maßnahme-Teilnahme-Grunddatei) und ISAAK (Instrumente Aktiver Arbeitsmarktpolitik). Eine Betrachtung ausgewählter Fälle am Beispiel der Förderung im Rahmen des ESF-BA-Programms. FDZ Methodenreport Nr. 08/2008.
- Heckman, James J. / Ichimura, Hidehiko / Todd, Petra E. (1998): Matching as an Econometric Evaluation Estimator: Evidence from Evaluating a Job Training Programme. In: *The Review of Economic Studies* 65 (24). 261-294.
- Imbens, Guido W. (2004). Nonparametric Estimation of Average Treatment Effects under Exogeneity: a Review. In: *The Review of Economics and Statistics* 86 (1). 4-29.
- Jacobebbinghaus, Peter / Seth, Stefan (2007): The German integrated employment biographies sample IEBS. In: *Schmollers Jahrbuch* 127 (2). 335-342.
- Jaenichen, Ursula / Kruppe, Thomas / Stephan, Gesine / Ullrich, Britta / Wießner, Frank (2005): You can split it if you really want: Korrekturvorschläge für ausgewählte Inkonsistenzen in IEB und MTG. FDZ Datenreport Nr. 04/2005.
- Kluve, Jochen (2006): The Effectiveness of European Active Labor Market Policy. IZA Discussion 2018
- Kruppe, Thomas / Oertel, Martina (2003): Von Verwaltungsdaten zu Forschungsdaten – Die Individualdaten für die Evaluation des ESF-BA-Programms 2000 bis 2006. IAB Werkstattbericht.
- Kruppe, Thomas / Müller, Eva / Wichert, Laura / Wilke, Ralf A. (2008): On the Definition of Unemployment and its Implementation in Register Data – The Case of Germany. In: *Schmollers Jahrbuch* 128 (3), 461-488.
- Nordberg Leif (2003): An Analysis of the Effects of using Interview versus Register Data in Income Distribution Analysis based on the Finnish ECHP-Surveys in 1996 and 2000. CHINTEX Working Paper nr. 15.
- Roed, Knut / Raaum, Oddbjørn (2003): Administrative Registers- Unexplored Reservoirs of Scientific Knowledge? In: *The Economic Journal* 113. 258-281.
- Rosenbaum, Paul R. / Rubin, Donald B. (1985): Constructing a Control Group Using Multivariate Matched Sampling Methods that Incorporate the Propensity Score. In: *The American Statistician* 39 (1). 33-38.
- Schnell, Rainer (1985): Zur Effizienz einiger Missing – Data -Techniken - Ergebnisse einer Computer – Simulation, In: *ZUMA-Nachrichten* 17. 50-74.
- Schnell, Rainer (1991): Der Einfluß gefälschter Interviews auf Survey-Ergebnisse. In: *Zeitschrift für Soziologie* 20 (1). 25-35.
- Waller, Marie (2007): Do Reported End dates of Treatments Matter for Evaluation Results? – An Investigation Based on the German Integrated Employment Biographies Sample. FDZ Methodenreport 01/2007.
- Wunsch, Conny / Lechner, Michael (2008): What did all the money do? On the General Ineffectiveness of Recent West German Labour Market Programmes. In: *Kyklos*, Blackwell Publishing 61 (1). 134-174.
- Zimmermann, Ralf / Kaimer, Steffen / Oberschachtsiek, Dirk (2007): Dokumentation des „Scientific Use Files der integrierten Erwerbsbiographien“ (IEBS-SUF V1) Version 1.0. FDZ Datenreport 01/2007 (de).

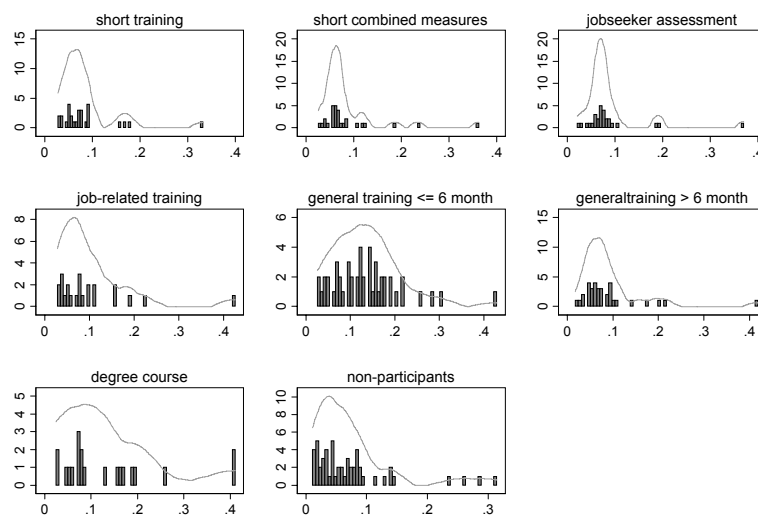
Appendix

Table A1: Variables

Personal Characteristics	
Name	definition
age (N)	age at the point in time of programme start
gender (D)	1=female
disabled (D)	1=individual is disabled
foreign (D)	1=individual is not German
health (D)	1=individual has health problems
health affection (D)	1=health problems affect employability
family status (N)	married, single, single with child, couple
education (N)	no professional degree, completed apprenticeship, Univer-
profession (N)	unskilled, technical profession, services-construction-
occupational status (N)	high-skilled, unskilled, clerk, part-time worker
desired occupation (N)	admin-teaching-science, other services, manufacturing-
wage (N)	last wage
employment status (N)	employed, unemployed, type of programme, type of em-
begin/end dates (N)	Begin and end dates of programme, employment, unem-
desired job (N)	unskilled, skilled, high-skilled, full-time only, part-time
ubclaim (N)	remaining unemployment benefit claim in days
sanction (D)	1=at least one benefit sanction
no attendance (D)	1=did not attend interview at PES at least once
no cooperation (D)	1=individual was not cooperative at least once
Statistics on regional level	
Name	definition
unemployment rate (N)	local unemployment rate
employment rate (N)	local employment rate
industry (N)	industry quota
household income (N)	household income per capita in euro
income tax (N)	income tax per capita in euro

Fortsetzung Tabelle A1	
Name	definition
GDP (N)	GDP per capita in euro
rural area (N)	rural population
federal state (N)	Schleswig-Holstein, Hamburg, Lower Saxony-Bremen,
medium (N)	number of medium-size cities
social assistance (N)	social assistance recipients per capita

Figure A1: Differences of significant dummy-variables



Source: IEBS, Wunsch/Lechner (2007), own calculations

Table A2: Means and Shares of Selected Variables

A) short training						
	method V0 (n=1126)		method V1 (n=1079)		method V2 (n=1113)	
Variable	mean	std. dev.	mean	std. dev.	mean	std. dev.
Age	36.51	6.632	36.56	6.665	36.51	6.614
female	0.51	0.500	0.50	0.500	0.51	0.500
married	0.47	0.499	0.47	0.499	0.48	0.500
completed apprenticeship	0.64	0.479	0.64	0.480	0.64	0.481
health problems	0.15	0.359	0.15	0.362	0.15	0.362
foreign	0.12	0.324	0.12	0.329	0.12	0.327
fulltime	0.76	0.430	0.77	0.420	0.77	0.423
clerk	0.39	0.487	0.38	0.486	0.39	0.488
claim in days	85.69	119.723	83.53	119.978	83.68	119.498
last monthly earnings	1.734	862.35	1.710	1017.43	1.723	993.88
duration last unemployment	4.94	3.661	4.75	3.662	4.88	3.680
total time unemployed	17.88	18.003	17.77	17.949	17.98	18.216
unemployment rate	8.83	2.757	8.87	2.767	8.85	2.736
employment rate	53.40	17.174	53.39	16.905	53.41	17.144
no significant differences (5%-level)						

B) short combined measures						
	method V0 (n=1366)		method V1 (n=1332)		method V2 (n=1268)	
Variable	mean	std. dev.	mean	std. dev.	mean	std. dev.
age	36.66	6.785	36.64	6.750	36.62	6.791
female	0.47	0.499	0.46	0.499	0.46	0.499
married	0.47	0.499	0.47	0.499	0.47	0.499
completed apprenticeship	0.57	0.495	0.57	0.496	0.57	0.495
health problems	0.20	0.402	0.20	0.402	0.21	0.406
foreign	0.14	0.352	0.14	0.351	0.14	0.346
fulltime	0.75	0.430	0.77	0.423	0.77	0.423
clerk	0.27	0.446	0.26	0.439	0.28	0.450
claim in days	66.45	107.643	64.82	105.941	63.77	105.554
last monthly earnings	1.605	813.54	1.558	844.59	1.605	817.06
duration last unemployment	5.38	3.801	5.08	3.713	5.29	3.760

<i>Fortsetzung Tabelle A2</i>						
total time unemployed	22.41	22.532	22.39	22.630	22.36	22.641
unemployment rate	8.58	2.856	8.62	2.851	8.61	2.862
employment rate	52.90	18.137	53.05	18.173	52.97	18.141
no significant differences (5%-level)						

C) Jobseeker assessment						
	method V0 (n=1529)		method V1 (n=1490)		method V2 (n=1513)	
Variable	mean	std. dev.	mean	std. dev.	mean	std. dev.
age	35.61	6.705	35.66	6.684	35.60	6.729
female	0.41	0.492	0.40	0.490	0.41	0.493
married	0.47	0.499	0.46	0.498	0.47	0.499
completed apprenticeship	0.55	0.497	0.55	0.498	0.54	0.498
health problems	0.19	0.396	0.19	0.393	0.20	0.399
foreign	0.10	0.305	0.10	0.299	0.11	0.310
fulltime	0.79	0.409	0.79	0.406	0.79	0.404
clerk	0.26	0.439	0.25	0.431	0.25	0.436
claim in days	98.52	124.339	95.60	123.894	95.95	124.032
last monthly earnings	1.638	796.90	1.586	856.53	1.619	809.96
duration last unemployment	5.18	3.965	4.84	3.925	5.12	3.977
total time unemployed	21.12	20.507	21.33	21.174	21.02	20.530
unemployment rate	9.35	2.728	9.35	2.736	9.34	2.739
employment rate	50.47	15.232	50.35	14.825	50.33	15.059
significant on a 5%-level: last earning < 1000						

D) job-related training						
	method V0 (n=603)		method V1 (n=582)		method V2 (n=594)	
Variable	mean	std. dev.	mean	std. dev.	mean	std. dev.
age	37.47	6.641	37.59	6.622	37.63	6.636
female	0.44	0.497	0.44	0.497	0.44	0.497
married	0.48	0.500	0.48	0.500	0.49	0.500
completed apprenticeship	0.63	0.484	0.62	0.486	0.63	0.483
health problems	0.14	0.345	0.14	0.343	0.15	0.354
foreign	0.13	0.336	0.13	0.331	0.13	0.332

<i>Fortsetzung Tabelle D</i>						
fulltime	0.74	0.438	0.75	0.434	0.76	0.427
clerk	0.26	0.441	0.26	0.441	0.27	0.444
claim in days	171.43	153.197	168.40	154.865	169.17	154.370
last monthly earnings	1.699	1007.50	1.649	1047.42	1.688	1024.35
duration last unemployment	5.12	3.939	4.82	3.847	5.06	3.969
total time unemployed	19.61	20.003	19.52	20.310	19.96	20.641
unemployment rate	8.67	2.758	8.66	2.738	8.64	2.727
employment rate	50.23	14.979	50.22	14.917	50.37	15.158
no significant differences (5%-level)						

E) degree course						
	method V0 (n=548)		method V1 (n=520)		method V2 (n=533)	
Variable	mean	std. dev.	mean	std. dev.	mean	std. dev.
age	33.67	6.065	33.75	6.051	33.60	6.061
female	0.45	0.498	0.45	0.498	0.45	0.498
married	0.45	0.498	0.44	0.497	0.45	0.498
completed apprenticeship	0.44	0.497	0.43	0.496	0.43	0.496
health problems	0.08	0.275	0.08	0.276	0.08	0.275
foreign	0.12	0.324	0.11	0.317	0.13	0.334
fulltime	0.48	0.500	0.49	0.500	0.48	0.500
clerk	0.23	0.421	0.21	0.409	0.23	0.419
claim in days	147.84	124.735	144.23	125.484	146.16	125.269
last monthly earnings	1.694	817.68	1.619	88.47	1.659	848.83
duration last unemployment	5.17	3.822	4.75	3.722	5.09	3.825
total time unemployed	17.50	16.756	17.51	17.263	16.96	16.318
unemployment rate	8.77	2.627	8.88	2.697	8.87	2.693
employment rate	52.52	17.088	52.98	17.114	52.27	16.668
no significant differences (5%-level)						

F) Non-participants						
	method V0 (n=22095)		method V1 (n=20222)		method V2 (n=21682)	
Variable	mean	std. dev.	mean	std. dev.	mean	std. dev.
age	36.53	6.836	36.57	6.839	36.56	6.831
female	0.49	0.500	0.48	0.500	0.49	0.500
married	0.52	0.499	0.51	0.500	0.52	0.500
completed apprenticeship	0.52	0.500	0.51	0.500	0.52	0.500
health problems	0.21	0.405	0.21	0.409	0.21	0.407
foreign	0.18	0.382	0.18	0.385	0.18	0.381
fulltime	0.75	0.435	0.76	0.429	0.75	0.434
clerk	0.23	0.419	0.23	0.418	0.23	0.419
claim in days	40.71	83.141	39.33	82.251	39.64	81.885
last monthly earnings	1.422	1049.49	1.406	1077.97	1.418	1060.08
duration last unemployment	4.97	2.149	4.71	2.179	4.90	2.147
total time unemployed	21.11	21.511	20.66	21.369	20.82	21.370
unemployment rate	9.00	2.760	9.02	2.767	9.01	2.763
employment rate	51.91	16.472	52.11	16.621	51.93	16.497
significant on a 5%-level: last profession. year of entry. duration last unemployment and employment. unemployed 6 months before programme (5%-level)						